

CLAIMS

1 1. A method for load balancing a plurality of servers, the method comprising:
2 providing a plurality of control blocks, each control block associated with a num-
3 ber of active connections a server is connected with, the control block configured to con-
4 trol at least one server with the associated number of connections in a server list;
5 causing each control block to point to a server with a least value ascertained by
6 determining the number of connections on the server relative to the server's capacity to
7 handle connections;
8 selecting the control block associated with the least number of connections; and
9 selecting the server pointed to by the control block.

1 2. The method as in claim 1, wherein ascertaining the least value for the
2 server comprises:
3 determining a metric of the server by dividing the number of connections on the
4 server by the capacity of the server, wherein the metric is kept as a quotient/remainder
5 pair;
6 storing the quotient/remainder pair in the control block;
7 incrementing the remainder by one for every connection added to the server; and
8 decrementing the remainder by one for every connection removed from the server.

1 3. The method as in claim 1, further comprising:
2 causing the control block with the server having an added/removed connection to
3 transfer the server to an adjacent control block, wherein the adjacent control block is as-
4 sociated with the number of connections pertaining to the transferring server;
5 causing the control block to transfer the metric of the server to the adjacent con-
6 trol block; and
7 updating the pointer to point to the next server on the list of the control block.

1 4. The method as in claim 3, further comprising:
2 removing the control block if the control block does not have a server on the server list.

1 5. The method as in claim 3, further comprising:
2 causing the adjacent control block to receive the transferring server;
3 causing the adjacent control block to receive the metric of the transferring server;
4 and
5 causing the adjacent control block to update and sort the server list.

1 6. The method as in claim 5, further comprising:
2 adding a control block if there is no control block associated with the number of connec-
3 tions of the transferring server.

1 7. A processor executable medium which when executed by a processor per-
2 forms a method for load balancing a plurality of servers, the method comprising:
3 providing a plurality of control blocks, each control block associated with a num-
4 ber of active connections a server is connected with, the control block configured to con-
5 trol at least one server with the associated number of connections in a server list;
6 causing each control block to point to a server with a least value ascertained by
7 determining the number of connections on the server relative to the server's capacity to
8 handle connections;
9 selecting the control block associated with the least number of connections; and
10 selecting the server pointed to by the control block.

1 8. The processor executable medium as in claim 7, wherein ascertaining the
2 least value for the server comprises:
3 determining a metric of the server by dividing the number of connections on the
4 server by the capacity of the server, wherein the metric is kept as a quotient/remainder
5 pair;
6 storing the quotient/remainder pair in the control block;
7 incrementing the remainder by one for every connection added to the server; and
8 decrementing the remainder by one for every connection removed from the server.

- 17